1  **Translational efficiency across healthy and tumor tissues is**

2  **proliferation-related**

3  **Xavier Hernandez-Alias[1], Hannah Benisty[1], Martin H. Schaefer[1,2,*], Luis Serrano[1,3,4,*]**

4  [1]Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology,

5  Dr. Aiguader 88, Barcelona 08003, Spain

6  [2]IEO European Institute of Oncology IRCCS, Department of Experimental Oncology, Via

7  Adamello 16, Milan 20139, Italy

8  [3]Universitat Pompeu Fabra (UPF), Barcelona 08002, Spain

9  [4]ICREA, Pg. Lluís Companys 23, Barcelona 08010, Spain

10  *Corresponding authors: martin.schaefer@ieo.it (M.H.S.), luis.serrano@crg.eu (L.S.)

11  **ABSTRACT**

12  **Background:** Different tissues express genes with particular codon usage and anticodon

13  tRNA repertoires. However, the codon-anticodon co-adaptation in humans is not completely

14  understood, as well as its effect on tissue-specific protein levels.

15  **Results:** We first validated the accuracy of small RNA-seq for tRNA quantification across

16  five human cell lines. We then analyzed tRNA expression in more than 8000 tumor samples

17  from TCGA, together with their paired mRNA-seq and proteomics data, to determine the

18  Relative Translation Efficiency. We thereby elucidate that the dynamic adaptation of the

19  tRNA pool is largely related to the proliferative state across tissues, which determines tissue-

20  specific translation efficiency. Furthermore, the aberrant translational efficiency of ProCCA

21  and GlyGGT in cancer, among other codons, which is partly regulated by the tRNA gene

22  copy numbers and their promoter DNA methylation, is associated with poor patient survival.

23  **Conclusions:** The distribution of tissue-specific tRNA pools over the whole cellular

24  translatome affects the subsequent translational efficiency, which functionally determines a

25  condition-specific expression program in tissues both in healthy and tumor states.

26    **KEYWORDS**

27    tRNA, translation, The Cancer Genome Atlas, tissue, regulation, miRNA, codon usage

28    **BACKGROUND**

29    In the light of the genetic code, multiple 3-letter combinations of nucleotides in the mRNA

30    can give rise to the same amino acid, which are known as synonymous codons. However,

31    despite the homology at the protein level, these different codons are recognized distinctly by

32    the transcriptional and translational machineries (1,2), and ultimately cause changes at

33    multiple levels of gene expression. Therefore, the non-uniform abundance of synonymous

34    codons across different tissues and among distinct functional gene sets has been proposed

35    as an adaptive mechanism of gene expression regulation (3), particularly linked to the

36    proliferative state (4). Nevertheless, in human, it is still under debate whether the efficiency

37    of gene expression is the main selective pressure driving the evolution of genomic codon

38    usage (5).

39    The 61 amino-acid-coding codons need to be recognized by 46 different tRNA isoacceptors

40    distributed across 428 Pol-III-transcribed tRNA genes (6), thus requiring wobble interactions

41    (non-Watson-Crick base pairing). This complexity of the tRNA repertoire is further enhanced

42    by an average of 11-13 base modifications per tRNA and all possible combinations thereof

43    (7). The underlying mechanisms regulating tRNA gene expression and modification are far

44    from resolved (8,9). However, it has been established that different conditions and tissues

45    showcase distinct tRNA abundances (4,10) and codon usages (3,11).

46    In order to understand such changes in codon-anticodon co-adaptation, orthogonal datasets

47    of gene expression including tRNA quantification are required, which needs to overcome the

48    challenges of strong secondary structures and abundant chemical modifications. Recent

49    technological developments have paved the way for sensitive high-throughput tRNA

50    sequencing across tissues and conditions (12,13). Aside from these methods and despite

51    the lower coverage, tRNA reads can also be detected from generic small RNA-seq datasets

52    (14–18). In this context, The Cancer Genome Atlas (TCGA) has been recently used to

53    investigate the alteration of tRNA gene expression and translational machinery in cancer,

54    which may play a role in driving aberrant translation (19,20).

55    To validate the use of small RNA-seq for tRNA quantification, we first compare tRNA levels

56    determined in HEK293 by well-established tRNA sequencing methods (Hydro-tRNAseq and

57    demethylase-tRNA-seq) (12,13,21), with those obtained by small RNA-seq. Then we

58    quantify the tRNA repertoire of five cell lines using Hydro-tRNAseq and perform small RNA-

59    seq in parallel. Comparison of the tRNA abundance obtained by both approaches shows that

60    it is possible to accurately estimate relative tRNA abundance of cells and tissues using small

61    RNA-seq. Furthermore, we show that both types of quantification are informative enough to

62    distinguish between the five analyzed human cell lines covering multiple tissue types. In

63    consequence, we apply a tRNA-specific computational pipeline to re-analyze 8,534 small

64    RNA-seq datasets from TCGA (22). We find that the tissue-specificity of tRNA expression is

65    largely proliferation-related, even within healthy tissues. The tRNA quantification of TCGA

66    samples enables their comparison with paired and publicly available mRNA-seq, proteomic,

67    DNA methylation and copy number data, which underscores the role of tRNAs in globally

68    controlling a condition-specific translational program. We discover multiple codons, including

69    ProCCA and GlyGGT, whose translational efficiency is compromised and leads to poor

70    prognosis in cancer. Finally, promoter DNA methylation and tRNA gene copy number arise

71    as two regulatory mechanisms controlling tRNA gene expression in cancer.

72    **RESULTS**

73    **tRNA quantification and modifications from small RNA-seq data**

74    In order to test how accurately we can extract tRNA abundance information contained in

75    small RNA sequencing data, we re-analyze four publicly-available datasets of the cell line

76    HEK293 (18,23,24). In contrast to previous studies analyzing tRNA expression from small

77    RNA-seq data (19,20), we use a computational pipeline specifically developed for the

3

78    accurate mapping of tRNA reads (16) in order to quantify all different isoacceptor species

79    (Figure 1A, see Methods). To validate the accuracy of these small RNA-seq quantifications,

80    we retrieve three datasets of well-established tRNA sequencing methods (Hydro-tRNAseq

81    and demethylase-tRNA-seq) applied to the same cell type (12,13,21), which autocorrelate in

82    the range of 0.72-0.79 among themselves (Table S1). In comparison, our four HEK293 small

83    RNA-seq quantifications show an average Spearman correlation against these three

84    conventional datasets of 0.68. Compared to the anticodon Spearman correlation of 0.61 as

85    published by Zhang et al. (19), our tRNA-specific mapping pipeline performs better than the

86    previously published protocol, since all our correlations lie over that value.

87    Further than correlating small RNA-seq data with conventional tRNA-seq datasets, we

88    analyze whether small RNA-seq quantifications are informative enough to distinguish

89    between different human cell lines covering multiple tissue types. We therefore apply both

90    small RNA-seq and Hydro-tRNAseq to HEK293 (kidney), HCT116 (colon), HeLa (cervix),

91    MDA-MB-231 (breast), and BJ fibroblasts. However, given the high variability between

92    replicates of MDA-MB-231 Hydro-tRNAseq quantifications, this cell line was excluded from

93    further analyses (Table S2). First, the correlations between the two methods of identical

94    samples and computational mapping pipeline range between 0.93 and 0.96 for all cell lines.

95    tRNA quantifications from both protocols are compared and significantly higher Spearman

96    correlations are obtained within matching samples versus mismatching cell lines (Figure 1B).

97    In consequence, we demonstrate that small RNA-seq quantifications of sample-specific

98    tRNA profiles show a good agreement with Hydro-tRNAseq.

99    We also detect tRNA base modifications in both protocols by nucleotide variant calling, as

100    described in Hoffmann et al. (16). In all cases, considering the modifications that are

101    detected in all three replicates, Hydro-tRNAseq datasets identify a larger number of

102    modifications than small RNA-seq, as expected by the more uniform and deeper coverage of

103    this method (Table S2). Furthermore, we detect a significant enrichment of the Hydro-

104    tRNAseq modifications in the small RNA-seq data ($p < 1e-16$, Fisher test), indicating that the

105    latter contains also information on tRNA modifications (Figure 1C).

106    Taken together, these observations demonstrate the applicability of small RNA-seq data for

107    the quantification tRNAs and their modifications. We therefore apply the same computational

108    pipeline to all healthy and primary tumor small RNA-seq samples from 23 cancer types of

109    The Cancer Genome Atlas (TCGA), which consists of 8,605 samples distributed among 17

110    different human tissues (Figure 1D, number of samples and their abbreviations in Table S3).

111    **Figure 1. tRNA quantification and modifications from small RNA-seq data.** (A)

112    Schematic pipeline for accurate mapping of tRNA reads. (B) Correlations between tRNA

113    quantifications by small RNA-seq and Hydro-tRNAseq of matching (correlations within the

114    same cell line) versus non-matching (different cell lines) samples. The p-value corresponds

115    to a one-tailed Wilcoxon rank-sum test, with $n_{matching} = 9$ and $n_{non-matching} = 72$. (C) Overlap of

116    the detected tRNA modifications upon variant calling by both methods. (D) The TCGA

117    network contains small RNA-seq data alongside mRNA-seq, DNA methylation arrays, non-

118    targeted proteomics, and copy number alteration quantification comprising 17 tissues.

119    **Proliferation is the major driver of tissue-specificity in tRNAs**

120    To determine the tissue-specificity of tRNAs in physiological conditions, the tRNA levels of

121    all 675 healthy samples in TCGA tissues are analyzed.  The isoacceptor abundances show

122    a significant tissue-specificity for all 46 annotated anticodons ($q<0.05$, FDR-corrected

123    Kruskal-Wallis test grouped by cancer types). Such differences between tissues are also

124    observed by hierarchical clustering of the median expression between all groups (Figure

125    2A). Furthermore, healthy samples from cancer types originating from the same tissue tend

126    to cluster together: READ and COAD from the gut; KIRC, KIRP and KICH from the kidney;

127    LUAD and LUSC from the lung; UCEC and CESC from the uterus; LIHC and CHOL from the

128    liver (refer to Table S3 for full cancer names). On the other hand, in terms of anticodon

129  abundances, three main subgroups of tRNAs with low, medium and high expression can be

130  distinguished across all cancer types (Figure 2A).

131  Regarding codon usage, a measure of tRNA abundance taking into account the relative

132  contribution of each tRNA anticodon among the set of codons of a certain amino acid is the

133  Relative Anticodon Abundance. From this perspective, a principal component analysis (PCA)

134  of the healthy control samples in TCGA also shows clear differences between tissues

135  (Figure 2B). Moreover, our first component, which explains 18.5% of the variance, correlates

136  positively with the proliferation marker Ki67 ($R_{spearman}$= 0.45) (25). To further interrogate the

137  biological functions related to the variability of anticodon abundances between samples, we

138  compute the correlation of the whole mRNA-seq transcriptome against the first PCA

139  component, and analyze it by Gene Set Enrichment Analysis (GSEA). As a result, the top

140  correlating genes are enriched in proliferation and immune cell activation, while the lowest

141  correlations belong to genes related with oxidative metabolism and respiration (Figure 2C,

142  Table S4). This confirms, as has been previously suggested (4), that there is a proliferative

143  tRNA expression program.

144  Overall, we observe patterns of tissue-specific tRNA expression in TCGA healthy samples.

145  Furthermore, our analysis identifies the proliferative state of tissues as the major biological

146  function driving the variability on tRNA abundances.

147  **Figure 2. Proliferation is the major driver of tissue-specificity in tRNAs.** (A) Medians of

148  square-root-normalized tRNA abundances across all TCGA tissues. The color of the tissue

149  labels correspond to the average Ki67 expression. (B) Principal Component Analysis (PCA)

150  of the Relative Anticodon Abundances (RAA) of TCGA, where the color scale corresponds to

151  the mean tissue expression of Ki67. The Spearman correlations of Ki67 with the components

152  are shown, as well as the samples of most extreme tissues. (C) Top positive and negative

153  GO terms upon Gene Set Enrichment Analysis (GSEA) of the correlations of the first PCA

154  component against all genes. Refer to Supplementary Table 3 for full cancer type names.

**tRNA repertoires determine tissue-specific translational efficiency**

155

156 Given that different tissues express distinct tRNA repertoires, we wondered whether they

157 could have an effect in protein translation. In this context, and based on previous studies

158 underscoring the global control role of codon usage as a competition for a limited tRNA pool

159 (26–28), we define the Relative Translation Efficiency (RTE) as the balance between the

160 supply (i.e. the anticodon tRNA abundances) and demand (i.e. the weighted codon usage

161 based on the mRNA levels) for each of the 60 codons (excluding methionine and Stop

162 codons). Furthermore, we normalize both the codon and anticodon abundances within each

163 amino acid family (i.e. relative to the most abundant synonymous codon/anticodon), in order

164 to remove the effect of amino acid biases and get a cleaner measure of codon optimality

165 (29).

166 To validate the suitability of RTE in determining the translational efficiency, we correlate the

167 RTE value of all proteins against the available proteomics data of paired TCGA samples

168 (30,31), which includes breast and colorectal tissues (tumor only, as no healthy samples are

169 available). Although the correlation is poor (but significant), both the protein abundances and

170 the protein-to-mRNA ratios correlate significantly better with RTE than with the classical

171 tRNA Adaptation Index [tAI] (32,33) or with a relative tAI with normalized weights within each

172 amino acid family, which do not consider the mRNA codon demand (Figure 3A).

173 Furthermore, the correlation of RTE with protein-to-mRNA ratio is slightly but significantly

174 higher than with protein levels alone, which indicates that the first is a better proxy for the

175 process of translation.

176 Next, we calculate the RTE for the 620 healthy samples for which both tRNA abundances

177 and mRNA levels are available. When analyzing the tissue medians of RTE weights per

178 each codon (RTEw), we observe that most codons are optimally balanced (RTEw =1), while

179 12.4% and 23.6% of codons are favored (RTEw >2) and disfavored (RTEw <0.5)

180 respectively. The tissue clustering again shows that healthy samples of cancer types from

7

181    the same tissue have similar RTEw profiles, which separates two major clusters of mostly

182    high-Ki67 and low-Ki67 tissues (Figure S1).

183    In order to identify the codons contributing most to the differences between tissues, we

184    compute a bidimensional PCA across all samples and RTEw (Figure 3B). Both the first and

185    second components significantly correlate with the proliferation marker Ki67 (0.4 and 0.35;

186    see Figure 2B). In agreement with the proliferation- and differentiation-related codons of

187    Gingold et al. (4), such proliferative pattern is similarly reproduced by the codons

188    contributing to the first PCA component, which has the strongest association to proliferation

189    (Figure 3B). Further, similarly to the tRNA abundances (Figure 2B), a GSEA of correlating

190    genes with the first component highlights the link with proliferation-related terms (Table S5).

191    On the other hand, the first component also clearly separates codons based on the GC

192    content of the third codon base, which has recently been associated with differentiation (high

193    in nnC/G codons) versus self-renewal functions (high in nnA/T) (34), as well as with

194    proliferative transcriptomes (35).

195    The previous analyses support the idea of proliferation-related tRNAs driving changes in

196    translational efficiencies. In that case, we expect that the two most extreme tissues in terms

197    of proliferation (brain and gut, excluding thymus for its low number of samples) differ in the

198    optimization of proliferation-related proteins. As such, we compute the average RTEw for

199    these two tissues, analyze the subsequent RTE score for each protein, and perform a GSEA

200    of the differential RTE per protein. Consistent with our hypothesis, the results indicate that

201    gut-optimized proteins are enriched in translation, DNA replication and protein localization,

202    whereas brain-optimized proteins are related to phospholipid production and neural function

203    (Figure 3C, Table S6). Taken together, this result confirms that the tRNA-dependent

204    translational efficiency is optimized for the translation of tissue-specific genes, particularly in

205    function of the proliferation state.

206 **Figure 3. tRNA repertoires determine tissue-specific translational efficiency.** (A) Three

207 metrics of translation efficiency (the classical tAI, a relative tAI with normalized weights

208 within each amino acid family, and the Relative Translation Efficiency described in this

209 article) are Spearman correlated against two proxies of translation (protein abundance and

210 protein-to-mRNA ratio) for all samples for which proteomics data is available (BRCA, COAD

211 and READ). Statistical differences are determined by sample-paired two-tailed Wilcoxon

212 rank-sum test. (B) Principal Component Analysis (PCA) of the RTEw of TCGA, where the

213 color scale corresponds to the mean tissue expression of Ki67. The Spearman correlations

214 of Ki67 with the components are shown, as well as the samples of most extreme tissues. On

215 the right, the top and bottom proliferation- and differentiation-related codons, as defined by

216 Gingold et al. (2014), ordered by their contribution to the first PCA component. (C) GSEA of

217 the differential RTE between extreme tissues ($\Delta RTE = RTE_{Colorectal} - RTE_{Brain}$), showing the

218 top five GO terms with high (left) and low (right) RTE in colorectal versus glial tissues. Refer

219 to Supplementary Table 3 for full cancer type names.

220 **Aberrant translational efficiencies drive tumor progression**

221 Given that proliferation is a major determinant of translational efficiency in healthy tissues, its

222 importance could be extrapolated to pathological conditions such as cancer. In fact, aberrant

223 expression of tRNAs and codon usage have been broadly related with tumorigenesis and

224 cancer progression (19,20,36,37). We therefore investigate 22 cancer types from TCGA in

225 order to determine which codons are translationally compromised in disease.

226 Similar to the analysis performed on the healthy tissues, we quantify all tRNA abundances of

227 TCGA primary tumor samples (Figure S2) and determine their corresponding translational

228 efficiencies using the RTE metric. By analyzing the differential RTEw between normal and

229 tumor samples, we observe many significant differences in all 60 codons across the 22

230 cancer types (Figure 4A). Among the most consistent changes, the ProCCA codon is

231 significantly more favored in tumors for 8 out of 10 cancer types, while the ProCCG is

9

232  disfavored in 14 out of 16 cancers (Figure 4B). In the case of glycine, translation appears

233  more efficient for GlyGGT in healthy samples (13/13), whereas tumor mostly favors GlyGGC

234  (9/12) and GlyGGG (7/9).

235  In terms of patient survival, we divide the TCGA patients in two groups based on their low or

236  high tumor RTEw and analyze their survival probability (Figure 4C, Table S7). Among

237  others, and consistent with the previous analysis, high translational efficiency weights of

238  ProCCA are associated with poor prognosis in kidney renal clear cell carcinoma and kidney

239  renal papillary cell carcinoma. Proline limitation in clear cell renal cell carcinoma has been

240  shown to compromise CCA-decoding tRNAPro aminoacylation, leading to reduced tumor

241  growth (38). In contrast, high RTEw of GlyGGT and ValGTC lead to longer survival in kidney

242  chromophobe and head and neck squamous cell carcinoma, respectively.

243  To determine the impact of aberrant translational efficiencies in regulating an oncogenic

244  translation program, we calculate the differential RTE for the whole genome based on the

245  average RTEw of healthy and tumor samples in kidney renal clear cell carcinoma, since it is

246  the cancer type with the most RTEw differences. The GSEA of the resulting ΔRTE score

247  indicates that cancer RTEw enhance the translation of proteins related to DNA replication

248  and gene expression, whereas the healthy kidney samples favor development and

249  differentiation processes (Table S8). As the RTEw of the ProCCA is specifically disturbed in

250  cancer, we also interrogate how this codon is distributed along the genome. We therefore

251  perform a GSEA on the relative codon usage of ProCCA, which shows that DNA replication

252  and cell cycle functions lie among the most CCA-enriched genes, while morphogenesis and

253  differentiation terms are CCA-depleted (Table S9). Together with the low-proliferative state

254  of kidney (Figure 2B), the over-efficiency of a proliferation-related codon in this tissue can

255  thus perturb its cellular RTE.

256  Overall, we detect differences at the level of RTEw between tumor and healthy tissues,

257  which show a functional relevance to the disease state. Therefore, while the differential

258    expression of tRNAs in TCGA had been already discussed elsewhere (19,20), we could here

259    elucidate their oncogenic effect in translational efficiency. In particular, ProCCA appears as

260    an interesting codon candidate in favoring tumor progression, which we had also detected in

261    healthy tissues to be associated with proliferation (Figure 3B, Table S5).

262    **Figure 4. Aberrant translational efficiencies drive tumor progression.** (A) Differential

263    RTEw between healthy and tumor samples across 22 cancer types, as measured by

264    $\log_2(\text{RTEw}_{Tumor}/\text{RTEw}_{Healthy})$. Only significant differences are colored, which are determined

265    using a two-tailed Wilcoxon rank-sum test and corrected for multiple testing by FDR. (B)

266    Boxplot of the RTEw of ProCCA and AlaGCG codons across TCGA cancer types. (C)

267    Survival curves for the previous codons in KIRC, KIRP and BLCA patients. The survival

268    analysis was performed for all codons whose translational efficiency was significantly

269    different in more than 5 cancer types in the one direction with respect to the other [Abs(UP-

270    DOWN)>5], and correspondingly corrected for multiple comparisons using FDR. Refer to

271    Supplementary Table 3 for full cancer type names.

272    **Promoter methylation and gene copy number regulate tRNA expression**

273    Aberrant translational efficiencies in cancer are partially caused by the differential expression

274    of tRNA genes (Figure S2). To determine the underlying mechanisms driving changes in

275    expression, we retrieve the methylation and copy number alteration (CNA) data from TCGA

276    samples, as a possible means for tRNA gene regulation. While CNA information cover 84%

277    of tRNA genes, the 450K-BeadChip methylation arrays used in TCGA are mostly centered

278    on the coding genome (Bibikova et al., 2011) and yield a coverage of only 37%.

279    In order to make the gene-based data comparable with the measured isoacceptor-based

280    tRNA expression, we average methylation and CNA levels over all genes within the same

281    isoacceptor family, at the cost of losing resolution. For each isoacceptor and each cancer

282    type, we finally fit a Multiple Linear Regression to determine how are promoter methylation

283    and CNA affecting tRNA expression (Figure 5A, Table S10). Among all models, the

11

284     significant coefficients for methylation and CNA are significantly negative and positive,

285     respectively. Despite the limited explained variance of the models (average $R^2$=0.023), such

286     results indicate that promoter methylation contributes to inhibition of tRNA gene expression,

287     while an increase in the gene copy number enhances tRNA expression.

288     Given the association of the codon ProCCA with cancer prognosis (Figure 4C), we explore

289     the expression pattern of tRNAPro in TCGA. In particular, tRNAPro$^{AGG}$, which recognizes the

290     codon ProCCA, is overexpressed in 8 out of 9 cancer types (Figure S2A). To get a more

291     accurate picture of the tRNA gene methylation levels, we also analyze recently published

292     bisulfite sequencing data (39), which, for 47 samples among nine cancer types, improved

293     the coverage of tRNA genes up to an average of 81%. In total, tRNAPro$^{AGG}$ genes stand

294     among the most duplicated and least methylated proline isoacceptors in cancer (Figure S3A-

295     B), in particular at the chr6.tRNA12 and chr16.tRNA12 genes (Figure 5B). Furthermore,

296     tRNAPro$^{AGG}$ gene duplications occur most frequently in kidney cancers (Figure S3C). On the

297     other hand, although the other CCA-decoding tRNAPro$^{TGG}$ is not differentially expressed in

298     cancer (Figure S2), its genes are as similarly methylated and duplicated as tRNAPro$^{AGG}$

299     (Figure 5B, Figure S3).

300     In short, promoter methylation and CNA appear as two possible regulatory mechanisms of

301     tRNA expression in cancer, which suggests that similar mechanisms that control the Pol-II-

302     mediated RNAs might also regulate the expression of Pol-III non-coding transcriptome, such

303     as tRNA genes. However, more accurate and high-throughput data on the methylation and

304     CNA of the non-coding genome together with gene-based tRNA quantifications are needed

305     to make stronger associations.

306     **Figure 5. Promoter methylation and gene copy number regulate tRNA expression.** (A)

307     A Multiple Linear Regression (MLR) between square-root-normalized tRNA expression and

308     the average promoter methylation (450K BeadChip array) and gene copy number at the

309     isoacceptor level. Among all MLRs for each isoacceptor and each cancer type separately,

310    the dots show the FDR-normalized significant coefficients based on their corresponding t-

311    statistic p-value, and red/blue show whether they are negative/positive respectively. The p-

312    value corresponds to a two-tailed binomial test between $n_{pos}$ and $n_{neg}$. (B) Differential

313    promoter methylation (bisulfite sequencing) between healthy and tumor samples of genes

314    expressing proline tRNAs, as measured by $\Delta\%Me=(\%Me_{Tumor}-\%Me_{Healthy})$. Refer to

315    Supplementary Table 3 for full cancer type names.

316    **DISCUSSION**

317    In this study, we use a systems biology approach to interrogate the multi-omics TCGA

318    dataset under the perspective of translational efficiencies. We therefore first validate the

319    suitability of small RNA-seq data in reproducing conventional tRNA-seq quantifications

320    based on a gold standard set of five tissue-wide human cell lines. In fact, knowing that small

321    RNA-seq datasets have a limited tRNA coverage and tend to be biased towards tRNA

322    fragments and unmodified tRNAs (18,40), we extend and apply a computational pipeline for

323    accurate mapping of tRNA reads (16). As a result, we obtain reproducible and informative

324    quantifications of all isoacceptors in our gold standard cell lines as well as in thousands of

325    samples across 23 cancer types of TCGA, exceeding the quality of similarly published data

326    (19,20).

327    From these quantifications, we then elucidate their effect on the translational efficiency by

328    defining the RTE, for Relative Translation Efficiency, which is a balance between the tRNA

329    supply and the codon demand. Although a more accurate RTE would have determined the

330    supply and demand based on the aminoacylated portion of tRNAs (41) and the ribosome-

331    bound mRNAs (42) respectively, we approximate such measures by our tRNA

332    quantifications and the publicly-available mRNA-seq data of TCGA. In agreement with

333    current studies showing that a dynamic codon usage need to compete for a limited tRNA

334    pool (28,29), we demonstrate that RTE is better measure of codon optimality than previously

335    published metrics such as the tAI (32,33). However, far from explaining the translation

336   process, the still low but significant correlations of protein-RTE in human, in contrast to

337   unicellular organisms, suggest that protein expression is also dependent on other layers of

338   regulation, such as transcriptional and post-transcriptional machineries, translation initiation,

339   epigenetic modifications of DNA and RNAs, or protein degradation mechanisms (43).

340   On the level of translational efficiency, in agreement with previous studies (4,36), we detect

341   that the proliferative state is the major determinant of RTE differences both across healthy

342   tissues and in cancer. Moreover, in contrast to recent work challenging the tissue-specificity

343   of codon-anticodon co-adaptation in human (29,43), our data here support the idea that

344   tissue-specific RTEw have functional implications on the tissue phenotype (e.g. in

345   determining neural differentiation in brain, or inducing abnormal proliferation in cancer).

346   Furthermore, we observe a pattern of proliferative nnA/T versus differentiative nnC/G

347   codons. Based on ribosome profiling experiments of pluripotency changes in embryonic

348   stem cells (34), this could be attributed to the slower translation in differentiated cells of

349   codons decoded by tRNAs that require adenosine-to-inosine modification at the wobble-

350   base pairing position. In particular, we detect the ProCCA codon to be significantly more

351   favored in proliferative cells and leading to poor cancer prognosis in kidney carcinomas,

352   specifically driven by an overexpression of tRNAPro$^{AGG}$ in cancer. Proline limitation in clear

353   renal cell carcinoma has indeed been shown to mostly compromise tRNAPro$^{AGG}$

354   aminoacylation, leading to slower proline translation and reduced tumor growth (38).

355   Furthermore, in support of our approach for isoacceptor quantification and translational

356   efficiency, similar studies of tRNA levels in TCGA have controversially claimed an opposite

357   prognostic value for the ProCCA codon in clear renal cell carcinoma (19,20).

358   In an effort to elucidate the mechanisms regulating the expression of tRNAs, we observe

359   that the tRNA gene copy number and their DNA methylation state have a positive and

360   inhibitory effect on tRNA expression, respectively. In this context, DNA methylation has

361   previously been linked to the silencing of type II genes (such as tRNAs) of the Pol-III

362   transcriptome (44). Here we specifically propose a role for DNA methylation in regulating the

14

363   overexpression of tRNAPro$^{AGG}$ in cancer. In terms of the copy number alterations, it is not

364   surprising to detect tRNA gene duplications in tumors, but the functional role in disease of

365   different isodecoder genes that share the same anticodon is still a matter of debate (45).

366   With the advent of more accurate and high-throughput multi-omics datasets, our knowledge

367   on the underlying mechanisms controlling tRNA expression, degradation, and the effect of

368   their modifications will be further expanded (8,9). Recent studies in TCGA have actually

369   observed an upregulation of tRNA-modifying enzymes, as well as proposed a link of tRNA-

370   derived fragments (tRF) to proliferation (19,46).

## CONCLUSIONS

371

372   This is the first high-throughput study of codon-anticodon translational efficiency over

373   thousands of samples comprising multiple tissues and disease. We therefore demonstrate a

374   functional role for the proliferation-driven tRNA expression differences in determining a

375   tissue-specific phenotype, both in physiological and pathological conditions. In the future, we

376   expect to validate the effect of such differential translational efficiency by integrating

377   perturbation☐based data and including additional gene expression regulatory layers such as

378   tRNA modifications.

## METHODS

379

### Cell lines

380

381   The cell lines included in this study are HeLa, HEK293, HCT116, MDA-MB-231 and

382   fibroblast BJ/hTERT. Cells were maintained at 37 °C in a humidified atmosphere at 5% $CO_2$

383   in DMEM 4.5g/L Glucose with UltraGlutamine media supplemented with 10% of FBS and 1%

384   penicillin/streptomycin.

### RNA extraction

385

386   Cells were grown in 60mm dishes for 48h. Total RNA from HeLa, HEK293, HCT116, MDA-

387   MB-231 and fibroblast BJ/hTERT was extracted using the miRNeasy Mini kit. Independent

388    replicates where grown and RNA was extracted on different days. 20 μg of total RNA was

389    treated following either the protocol of Hydro-tRNAseq (12) or generic small RNA-seq.

390    **Hydro-tRNA sequencing**

391    Total RNA was resolved on 15% Novex TBE urea gels and size-selected for 60-100 nt

392    fragments. The recovered material was then alkaline hydrolyzed (10mM sodium carbonate

393    and 10mM sodium bicarbonate) for 10 minutes at 60ºC. The resulting RNA was de-

394    phosphorylated with Antarctic Phosphatase (New England Biolabs) at 37ºC for 1 hour. De-

395    phosphorylated RNA was purified with an RNeasy MinElute spin column and re-

396    phosphorylated with Polynucleotide Kinase (NEB). PNK-treated tRNAs were purified with an

397    RNeasy MinElute spin column and, similar to small RNA-seq library preparation, adaptor-

398    ligated, reverse-transcribed and PCR-amplified for 14 cycles. The resulting cDNA was

399    purified using a QIAQuick PCR Purification Kit and sequenced on Illumina HiSeq 2500

400    platform in 50bp paired-end format.

401    From all five cell lines, the isoacceptor abundances of MDA-MB-231 yielded a median of 3-5

402    times higher standard deviation than the other Hydro-tRNAseq quantifications (Table S2),

403    thus suggesting some technical problem with this cell line. In consequence, this cell line was

404    excluded from any further analysis.

405    **Small RNA sequencing**

406    Total RNA was directly adaptor-ligated, reverse-transcribed and PCR-amplified for 12

407    cycles. The resulting cDNA was purified using a QIAQuick PCR Purification Kit and

408    sequenced on Illumina HiSeq 2500 platform in 50bp single-end format.

**409   Computational Analysis**

410   <u>tRNA quantification and modification calling</u>

411   In both Hydro-tRNAseq and small RNA-seq FASTQ files, sequencing adapters were
412   trimmed using BBDuk from the BBMap toolkit [v38.22]
413   (https://sourceforge.net/projects/bbmap): k-mer=10 (allowing 8 at the end of the read),
414   Hamming distance=1, length=10-50bp, Phred>25. Using the human reference genome
415   GRCh38 (Genome Reference Consortium Human Reference 38, GCA_000001405.15), a
416   total of 856 nuclear tRNAs and 21 mitochondrial tRNAs were annotated with tRNAscan-SE
417   [v2.0] (47).

418   Trimmed FASTQ files were then mapped using a specific pipeline for tRNAs (Figure 1A)
419   (16). Summarizing, an artificial genome is first generated by masking all annotated tRNA
420   genes and adding pre-tRNAs (i.e. tRNA genes with 3' and 5' genomic flanking regions) as
421   extra chromosomes. Upon mapping to this artificial genome with Segemehl [v0.3.1] (48),
422   reads that map to the tRNA-masked chromosomes or to the tRNA flanking regions are
423   filtered out in order to remove non-tRNA reads and unmature-tRNA reads respectively.

424   After this first mapping step, a second library is generated by adding 3' CCA tails and
425   removing introns from tRNA genes. All 100% identical sequences of this so-called *mature*
426   tRNAs are clustered to avoid redundancy. Next, the subset of filtered reads from the first
427   mapping is aligned against the clustered mature tRNAs using Segemehl [v0.3.1] (48).
428   Mapped reads are then realigned with GATK IndelRealigner [v3.8] (49) to reduce the
429   number of mismatching bases across all reads.

430   For quantification, isoacceptors were quantified as reads per million (RPM). In order to
431   increase the coverage for anticodon-level quantification, we consider all reads that map
432   unambiguously to a certain isoacceptor, even though they ambiguously map to different
433   isodecoders (i.e. tRNA genes that differ in their sequence but share the same anticodon).
434   Ambiguous reads mapping to genes of different isoacceptors were discarded.

17

435  Regarding modification site calling, we only considered gene-level uniquely mapped reads,

436  as described to be optimal in Hoffmann et al. (16). As in their pipeline, in order to distinguish

437  mapping or sequencing errors from true misincorporation sites, we use GATK

438  UnifiedGenotyper [v3.8] (49).

439  <u>Relative Codon Usage (RCU) and Relative Anticodon Abundance (RAA)</u>

440  The RCU/RAA is defined as the contribution of a certain codon/anticodon to the amino acid it

441  belongs to. The RCU of all synonymous codons and the RAA of all anticodons recognizing

442  synonymous codons therefore sum up to 1.

443
$$RCU = \frac{x_C}{\sum_{i \in C_{aa}} x_i} \qquad\qquad RAA = \frac{x_A}{\sum_{i \in A_{aa}} x_i}$$

444  where $x_C/x_A$ refers to the abundance of the codon/anticodon $C/A$, and $C_{aa}$ is the set of all

445  synonymous codons, as well as $A_{aa}$ is the set of all anticodons that decode synonymous

446  codons.

447  <u>tRNA Adaptation Index (tAI)</u>

448  As described by dos Reis et al. (2003, 2004), the tAI weights every codon based on the

449  wobble-base codon-anticodon interaction rules. Let $c$ be a codon, then the decoding weight

450  is a weighted sum of the square-root-normalized tRNA abundances $tRNA_{cj}$ for all tRNA

451  isoacceptors $j$ that bind with affinity $(1 - s_{cj})$ given the wobble-base pairing rules $n_c$.

452  However, while dos Reis et al. (2004) assumes that highly expressed genes are codon-

453  optimized, here we use the non-optimized s-values to avoid a circularity in our reasoning:

$$s = [0, 0, 0, 0, 0.5, 0.5, 0.75, 0.5, 0.5]$$

$$w_c = \sum_{j=1}^{n_c} (1 - s_{cj}) tRNA_{cj}$$

18

454    And therefore the tAI of a certain protein is the product of weights of each codon $i_k$ at the

455    triplet position $k$ throughout the full gene length $l_g$, and normalized by the length.

$$tAI = (\prod_{k=1}^{l_g} w_{i_k})^{1/l_g}$$

456    Relative tRNA Adaptation Index (RtAI)

457    For comparison with the RTE (Figure 3A), an amino-acid-normalized tAI measure is defined

458    by dividing each tAI weight by the maximum weight among all codons within each amino

459    acid family.

$$Rw_c = \frac{w_c}{max_{i \in c_{aa}}(w_i)}$$

460    And therefore the RtAI of a certain protein is the product of weights $Rw$ of each codon $i_k$ at

461    the triplet position $k$ throughout the full gene length $l_g$, and normalized by the length.

$$RtAI = (\prod_{k=1}^{l_g} Rw_{i_k})^{1/l_g}$$

462    Relative Translation Efficiency (RTE)

463    The RTE aims to consider not only tRNA abundances, but also the codon usage demand. In

464    doing so, it constitutes a global measure of translation control, since the efficiency of a

465    certain codon depends both on its complementary anticodon abundance as well as the

466    demand for such anticodon by other transcripts. This global control has been indeed

467    established to play an important role in defining optimal translation programs (28).

468    The definition of the RTE is based on similar previously published metrics (26,27), which

469    consists of a ratio between the anticodon supply and demand. On the one hand, the

470    anticodon supply is defined as the relative tAI weights $Rw$ (see previous section). On the

471    other, the anticodon demand is estimated from the codon usage at the transcriptome level. It

19

472    is computed as the frequency of each codon in a transcript weighted by the corresponding

473    transcript expression, and finally summing up over all transcripts. Let $c$ be a codon, then the

474    codon usage is a weighted sum of the counts of codon $c_i$ in gene $j$ weighted by the mRNA-

475    seq abundance $mRNA_j$ for all genes in the genome $g$:

$$CU_c = \sum_{j=1}^{g} c_{ij} mRNA_j$$

476    Similarly to the supply, the anticodon demand is then normalized within each amino acid

477    family:

$$D_c = \frac{CU_c}{max_{i \in c_{aa}}(CU_i)}$$

478    Finally, the RTE weights (RTEw) are defined as the ratio between the codon supply $S_c$ and

479    demand $D_c$:

$$RTEw_c = \frac{S_c}{D_c}$$

480    And therefore the RTE of a certain protein is the product of weights $RTEw$ of each codon $i_k$

481    at the triplet position $k$ throughout the full gene length $l_g$, and normalized by the length.

$$RTE = (\prod_{k=1}^{l_g} RTEw_{i_k})^{1/l_g}$$

482    Gene Set Enrichment Analysis (GSEA)

483    We analyzed the enrichment of gene sets of the GO Biological Process Ontology using the

484    GSEA algorithm (50). The score used to generate the ranked list input is specified in the text

485    for each analysis.

486    Survival Analysis

487    To analyze how translational efficiency of a certain codon (RTEw) can affect the survival

488    probability in cancer, patients of a certain cancer type are divided in two groups of low/high

489    RTEw, which correspond to the patients having the top and bottom 40% RTEw. The Kaplan-

490    Meier curves are then computed to estimate the survival probability of each group along

491    time.

492    tRNA methylation and copy number

493    For consistency with the current version of publicly available and pre-processed 450k DNA

494    methylation and SNP6 segmented copy number alteration (CNA) data from firebrowse, we

495    used the human reference genome GRCh37/hg19 (Genome Reference Consortium Human

496    Reference 37, GCA_000001405.1) in this analysis. The coordinates of all nuclear tRNA

497    genes were obtained using tRNAscan-SE [v2.0] (47).

498    Regarding DNA methylation, we computed the average beta value of each tRNA gene from

499    1.5kb upstream of the transcription start site (1500TSS) until the end of the gene. For CNA,

500    we retrieved the segmented data of precomputed $log_2(CN) - 1$ from firebrowse and

501    extracted the corresponding value for the genomic coordinates containing the tRNA genes.

502    Whenever the tRNA genes was located between two segments, the weighted average in

503    function of the gene overlap with each segment was computed.

504    Bisulfite sequencing methylation

505    As 1500TSS methylation of tRNA genes lead to an average coverage of only 37% genes,

506    we also analyzed the recently published bisulfite sequencing data of 47 samples across nine

507    cancer types (Table S3) (39). After retrieving the datasets from the GDC legacy archive,

508    given the higher resolution of bisulfite sequencing data, we restricted the computation of the

509    average promoter methylation of tRNA genes to the GRCh37/hg19 genomic coordinates

510    containing the tRNA genes, since the promoter region of Pol-III-genes is intragenic.

511    Multiple Linear Regression (MLR)

512    We fitted a Multiple Linear Regression (MLR) between the square-root-normalized tRNA

513    expression (dependent variable) and the promoter methylation and gene copy number

514    (independent variables). To make all three layers of information comparable, we considered

515    only samples for which all data was available and performed the regression at the

516    isoacceptor level, thus averaging the methylation and CNA data over all tRNA genes that

517    shared the same anticodon.

$$EXP = \beta_0 + \beta_{Me}Me + \beta_{CNA}CNA$$

518    We fitted the model parameters for all 64 isoacceptors and 22 cancer types, leading to

519    22x64=1408 MLRs, among which only significant coefficients (FDR-corrected t-statistic p-

520    value < 0.05) were considered in downstream analyses.

521    **Statistical Analysis**

522    For hypothesis testing, an unpaired two-tailed Wilcoxon rank-sum test was performed,

523    unless stated otherwise. All details of the statistical analyses can be found in the Results

524    section. We used a significance value of 0.05. In differential expression analyses, a False

525    Discovery Rate correction was used to account for multiple testing.

526    **ABBREVIATIONS**

527    TCGA: The Cancer Genome Atlas

528    PCA: Principal Component Analysis

529    GSEA: Gene Set Enrichment Analysis

530    RAA: Relative Anticodon Abundance

531    RTE: Relative Translation Efficiency

532    RTEw: RTE weights

533    tAI: tRNA Adaptation Index

534    RtAI: Relative tAI

22

535     CNA: Copy Number Alteration

536     BLCA: Bladder Urothelial Carcinoma

537     BRCA: Breast invasive carcinoma

538     CESC: Cervical squamous cell carcinoma and endocervical adenocarcinoma

539     CHOL: Cholangiocarcinoma

540     COAD: Colon adenocarcinoma

541     ESCA: Esophageal carcinoma

542     GBM: Glioblastoma multiforme

543     HNSC: Head and Neck squamous cell carcinoma

544     KICH: Kidney Chromophobe

545     KIRC: Kidney renal clear cell carcinoma

546     KIRP: Kidney renal papillary cell carcinoma

547     LIHC: Liver hepatocellular carcinoma

548     LUAD: Lung adenocarcinoma

549     LUSC: Lung squamous cell carcinoma

550     PAAD: Pancreatic adenocarcinoma

551     PCPG: Pheochromocytoma and Paraganglioma

552     PRAD: Prostate adenocarcinoma

553     READ: Rectum adenocarcinoma

554     SKCM: Skin Cutaneous Melanoma

555     STAD: Stomach adenocarcinoma

556     THCA: Thyroid carcinoma

557     THYM: Thymoma

558     UCEC: Uterine Corpus Endometrial Carcinoma

559 **DECLARATIONS**

560 **ETHICS APPROVAL AND CONSENT TO PARTICIPATE**

561 Not applicable.

562 **CONSENT FOR PUBLICATION**

563 Not applicable.

564 **AVAILABILITY OF DATA AND MATERIALS**

565 HEK293 Small RNA-seq datasets

566 The four HEK293 datasets were downloaded from the NCBI Sequence Read Archive (SRA):

567 SRR1304304, ERR705692, ERR705691, SRR2060090.

568 The Cancer Genome Atlas

569 Raw small RNA-sequencing data in BAM format were retrieved from the GDC legacy archive

570 after obtaining the necessary permissions from dbGaP, comprising all healthy samples (NT,

571 solid tissue normal) and their primary tumor (PT) counterparts, which consists of 23 cancer

572 types (BRCA, PRAD, KICH, KIRP, KIRC, LUAD, LUSC, HNSC, UCEC, CESC, LIHC, CHOL,

573 THCA, COAD, READ, ESCA, STAD, BLCA, PAAD, THYM, SKCM, PCPG, GBM). For

574 samples for which more than one BAM was available, all files were downloaded. BAM files

575 were converted to FASTQ using SAMtools [v1.3.1] (51). We retrieved publicly available and

576 pre-processed mRNA-seq gene expression, 450k DNA methylation, and SNP6 segmented

577 copy number alteration (CNA) from firebrowse. As for proteomics, preprocessed protein

578 assembly data and protein relative abundance were obtained from CPTAC for TCGA

579 samples including BRCA, COAD and READ.

580     Coding sequences

581     The coding sequences of *Homo sapiens* from RefSeq were downloaded from the

582     Codon/Codon Pair Usage Tables (CoCoPUTs) project release as of February 6, 2019

583     (52,53).

584     GO gene sets

585     Gene sets derived from the GO Biological Process Ontology were downloaded from the

586     Molecular Signatures Database [v6.2] (MSigDB) as a GMT file (50,54).

587     Generated data and code

588     The code used in this study is available at GitHub [https://github.com/hexavier/tRNA_TCGA;

589     https://github.com/hexavier/tRNA_mapping], and the generated datasets are publicly

590     accessible at Synapse (www.synapse.org/tRNA_TCGA, syn20640275). Hydro-tRNA and

591     small RNA sequencing data of all five cell lines has been made available at the Gene

592     Expression Omnibus (GEO): GSE137834. Hydro-tRNAseq data from HEK293 and HeLa has

593     been previously published (36) and deposited at ArrayExpress under accession number E-

594     MTAB-8144.

595     **COMPETING INTERESTS**

596     The authors declare that they have no competing interests.

597     **FUNDING**

603 **AUTHORS' CONTRIBUTIONS**

604 Conceptualization, X.H., M.H.S. and L.S.; Methodology, X.H., H.B., M.H.S., L.S.; Software,

605 X.H.; Investigation, H.B., X.H.; Validation, X.H., M.H.S.; Formal analysis, X.H., M.H.S.;

606 Writing-Original Draft, X.H.; Writing-Review & Editing, X.H., H.B., M.H.S., L.S.; Visualisation:

607 X.H., M.H.S.; Funding Acquisition, L.S.; Supervision, M.H.S. and L.S.

608 **ACKNOWLEDGEMENTS**

613 **REFERENCES**

614 1.  Hanson G, Coller J. Codon optimality, bias and usage in translation and mRNA decay.
615     Nat Rev Mol Cell Biol. 2017 Oct 11;19(1):20–30.
616 2.  Supek F. The Code of Silence: Widespread Associations Between Synonymous Codon
617     Biases and Gene Function. J Mol Evol. 2016 Jan;82(1):65–73.
618 3.  Najafabadi HS, Goodarzi H, Salavati R. Universal function-specificity of codon usage.
619     Nucleic Acids Res. 2009 Nov;37(21):7014–23.
620 4.  Gingold H, Tehler D, Christoffersen NR, Nielsen MM, Asmar F, Kooistra SM, et al. A
621     Dual Program for Translation Regulation in Cellular Proliferation and Differentiation.
622     Cell. 2014 Sep;158(6):1281–92.
623 5.  Pouyet F, Mouchiroud D, Duret L, Sémon M. Recombination, meiotic expression and
624     human codon usage. eLife [Internet]. 2017 Aug 15 [cited 2018 Nov 7];6. Available from:
625     https://elifesciences.org/articles/27344
626 6.  Chan PP, Lowe TM. GtRNAdb 2.0: an expanded database of transfer RNA genes
627     identified in complete and draft genomes. Nucleic Acids Res. 2016 Jan 4;44(D1):D184-
628     189.
629 7.  Schimmel P. The emerging complexity of the tRNA world: mammalian tRNAs beyond
630     protein synthesis. Nat Rev Mol Cell Biol. 2018 Jan;19(1):45–58.
631 8.  Pan T. Modifications and functional genomics of human transfer RNA. Cell Res. 2018
632     Apr;28(4):395.
633 9.  Rak R, Dahan O, Pilpel Y. Repertoires of tRNAs: The Couplers of Genomics and
634     Proteomics. Annu Rev Cell Dev Biol. 2018 Oct 6;34(1):239–64.
635 10. Dittmar KA, Goodenbour JM, Pan T. Tissue-Specific Differences in Human Transfer
636     RNA Expression. PLOS Genet. 2006 Dec 22;2(12):e221.
637 11. Waldman YY, Tuller T, Shlomi T, Sharan R, Ruppin E. Translation efficiency in
638     humans: tissue specificity, global optimization and differences between developmental
639     stages. Nucleic Acids Res. 2010 May 1;38(9):2964–74.
640 12. Gogakos T, Brown M, Garzia A, Meyer C, Hafner M, Tuschl T. Characterizing
641     Expression and Processing of Precursor and Mature Human tRNAs by Hydro-tRNAseq
642     and PAR-CLIP. Cell Rep. 2017 Aug;20(6):1463–75.

643   13.  Zheng G, Qin Y, Clark WC, Dai Q, Yi C, He C, et al. Efficient and quantitative high-
644        throughput tRNA sequencing. Nat Methods. 2015 Sep;12(9):835–7.
645   14.  Guo Y, Bosompem A, Mohan S, Erdogan B, Ye F, Vickers KC, et al. Transfer RNA
646        detection by small RNA deep sequencing and disease association with myelodysplastic
647        syndromes. BMC Genomics. 2015 Sep 24;16(1):727.
648   15.  Guo Y, Xiong Y, Sheng Q, Zhao S, Wattacheril J, Flynn CR. A micro-RNA expression-
649        signature for human NAFLD progression. J Gastroenterol. 2016 Oct;51(10):1022–30.
650   16.  Hoffmann A, Fallmann J, Vilardo E, Mörl M, Stadler PF, Amman F. Accurate mapping
651        of tRNA reads. Bioinformatics. 2018 Apr 1;34(7):1116–24.
652   17.  Pundhir S, Gorodkin J. Differential and coherent processing patterns from small RNAs.
653        Sci Rep [Internet]. 2015 Jul 13 [cited 2019 Jun 21];5. Available from:
654        https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4499813/
655   18.  Torres AG, Piñeyro D, Rodríguez-Escribà M, Camacho N, Reina O, Saint-Léger A, et
656        al. Inosine modifications in human tRNAs are incorporated at the precursor tRNA level.
657        Nucleic Acids Res. 2015 May 26;43(10):5145–57.
658   19.  Zhang Z, Ye Y, Gong J, Ruan H, Liu C-J, Xiang Y, et al. Global analysis of tRNA and
659        translation factor expression reveals a dynamic landscape of translational regulation in
660        human cancers. Commun Biol. 2018 Dec 21;1(1):234.
661   20.  Zhang Z, Ruan H, Liu C-J, Ye Y, Gong J, Diao L, et al. tRic: a user-friendly data portal
662        to explore the expression landscape of tRNAs in human cancers. RNA Biol. 2019 Aug
663        25;
664   21.  Mattijssen S, Arimbasseri AG, Iben JR, Gaidamakov S, Lee J, Hafner M, et al. LARP4
665        mRNA codon-tRNA match contributes to LARP4 activity for ribosomal protein mRNA
666        poly(A) tail length protection. eLife [Internet]. 2017 Sep 12 [cited 2019 Feb 27];6.
667        Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5626478/
668   22.  Chu A, Robertson G, Brooks D, Mungall AJ, Birol I, Coope R, et al. Large-scale
669        profiling of microRNAs for The Cancer Genome Atlas. Nucleic Acids Res. 2016 Jan
670        8;44(1):e3.
671   23.  Flores O, Kennedy EM, Skalsky RL, Cullen BR. Differential RISC association of
672        endogenous human microRNAs predicts their inhibitory potential. Nucleic Acids Res.
673        2014 Apr;42(7):4629–39.
674   24.  Mefferd AL, Kornepati AVR, Bogerd HP, Kennedy EM, Cullen BR. Expression of
675        CRISPR/Cas single guide RNAs using small tRNA promoters. RNA N Y N. 2015
676        Sep;21(9):1683–9.
677   25.  Scholzen T, Gerdes J. The Ki-67 protein: From the known and the unknown. J Cell
678        Physiol. 2000;182(3):311–22.
679   26.  Gingold H, Dahan O, Pilpel Y. Dynamic changes in translational efficiency are deduced
680        from codon usage of the transcriptome. Nucleic Acids Res. 2012 Nov 1;40(20):10053–
681        63.
682   27.  Pechmann S, Frydman J. Evolutionary conservation of codon optimality reveals hidden
683        signatures of cotranslational folding. Nat Struct Mol Biol. 2013 Feb;20(2):237–43.
684   28.  Frumkin I, Lajoie MJ, Gregg CJ, Hornung G, Church GM, Pilpel Y. Codon usage of
685        highly expressed genes affects proteome-wide translation efficiency. Proc Natl Acad
686        Sci U S A. 2018 22;115(21):E4940–9.
687   29.  Eraslan B, Wang D, Gusic M, Prokisch H, Hallström BM, Uhlén M, et al. Quantification
688        and discovery of sequence determinants of protein per mRNA amount in 29 human
689        tissues. Mol Syst Biol. 2019 Feb 1;15(2):e8513.
690   30.  Mertins P, Mani DR, Ruggles KV, Gillette MA, Clauser KR, Wang P, et al.
691        Proteogenomics connects somatic mutations to signalling in breast cancer. Nature.
692        2016 Jun;534(7605):55–62.
693   31.  Slebos RJC, Wang X, Wang X, Zhang B, Tabb DL, Liebler DC. Proteomic analysis of
694        colon and rectal carcinoma using standard and customized databases. Sci Data. 2015
695        Jun 23;2:150022.
696   32.  dos Reis M, Wernisch L, Savva R. Unexpected correlations between gene expression
697        and codon usage bias from microarray data for the whole Escherichia coli K-12

698        genome. Nucleic Acids Res. 2003 Dec 1;31(23):6976–85.

699  33.  dos Reis M, Savva R, Wernisch L. Solving the riddle of codon usage preferences: a
700        test for translational selection. Nucleic Acids Res. 2004 Jan 1;32(17):5036–44.

701  34.  Bornelöv S, Selmi T, Flad S, Dietmann S, Frye M. Codon usage optimization in
702        pluripotent embryonic stem cells. Genome Biol. 2019 Jun 7;20(1):119.

703  35.  Fornasiero EF, Rizzoli SO. Pathological changes are associated with shifts in the
704        employment of synonymous codons at the transcriptome level. BMC Genomics. 2019
705        Jul 9;20(1):566.

706  36.  Benisty H, Weber M, Hernandez-Alias X, Schaefer MH, Serrano L. Proliferation specific
707        codon usage facilitates oncogene translation. Manuscr Submitt Publ [Internet]. 2019 Jul
708        11 [cited 2019 Jul 24]; Available from:
709        https://www.biorxiv.org/content/10.1101/695957v1

710  37.  Goodarzi H, Nguyen HCB, Zhang S, Dill BD, Molina H, Tavazoie SF. Modulated
711        Expression of Specific tRNAs Drives Gene Expression and Cancer Progression. Cell.
712        2016 Jun;165(6):1416–27.

713  38.  Loayza-Puch F, Rooijers K, Buil LCM, Zijlstra J, F. Oude Vrielink J, Lopes R, et al.
714        Tumour-specific proline vulnerability uncovered by differential ribosome codon reading.
715        Nature. 2016 Feb;530(7591):490–4.

716  39.  Zhou W, Dinh HQ, Ramjan Z, Weisenberger DJ, Nicolet CM, Shen H, et al. DNA
717        methylation loss in late-replicating domains is linked to mitotic cell division. Nat Genet.
718        2018 Apr;50(4):591–602.

719  40.  Torres AG, Reina O, Attolini CS-O, Pouplana LR de. Differential expression of human
720        tRNA genes drives the abundance of tRNA-derived fragments. Proc Natl Acad Sci.
721        2019 Apr 23;116(17):8451–6.

722  41.  Evans ME, Clark WC, Zheng G, Pan T. Determination of tRNA aminoacylation levels by
723        high-throughput sequencing. Nucleic Acids Res. 2017 Aug 21;45(14):e133.

724  42.  Ingolia NT, Ghaemmaghami S, Newman JRS, Weissman JS. Genome-Wide Analysis
725        in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling. Science.
726        2009 Apr 10;324(5924):218–23.

727  43.  Rudolph KLM, Schmitt BM, Villar D, White RJ, Marioni JC, Kutter C, et al. Codon-
728        Driven Translational Efficiency Is Stable across Diverse Mammalian Cell States. Galtier
729        N, editor. PLOS Genet. 2016 May 11;12(5):e1006024.

730  44.  Park J-L, Lee Y-S, Kunkeaw N, Kim S-Y, Kim I-H, Lee YS. Epigenetic regulation of
731        noncoding RNA transcription by mammalian RNA polymerase III. Epigenomics. 2017
732        Jan 23;9(2):171–87.

733  45.  Lant JT, Berg MD, Heinemann IU, Brandl CJ, O'Donoghue P. Pathways to disease
734        from natural variations in human cytoplasmic tRNAs. J Biol Chem. 2019 Apr
735        5;294(14):5294–308.

736  46.  Telonis AG, Loher P, Magee R, Pliatsika V, Londin E, Kirino Y, et al. tRNA Fragments
737        Show Intertwining with mRNAs of Specific Repeat Content and Have Links to
738        Disparities. Cancer Res. 2019 Jan 1;canres.0789.2019.

739  47.  Chan PP, Lowe TM. tRNAscan-SE: Searching for tRNA Genes in Genomic Sequences.
740        Methods Mol Biol Clifton NJ. 2019;1962:1–14.

741  48.  Hoffmann S, Otto C, Kurtz S, Sharma CM, Khaitovich P, Vogel J, et al. Fast mapping of
742        short sequences with mismatches, insertions and deletions using index structures.
743        PLoS Comput Biol. 2009 Sep;5(9):e1000502.

744  49.  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The
745        Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
746        sequencing data. Genome Res. 2010 Sep;20(9):1297–303.

747  50.  Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al.
748        Gene set enrichment analysis: A knowledge-based approach for interpreting genome-
749        wide expression profiles. Proc Natl Acad Sci U S A. 2005 Oct 25;102(43):15545–50.

750  51.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
751        Alignment/Map format and SAMtools. Bioinforma Oxf Engl. 2009 Aug 15;25(16):2078–
752        9.

753   52.   Alexaki A, Kames J, Holcomb DD, Athey J, Santana-Quintero LV, Lam PVN, et al.
754         Codon and Codon-Pair Usage Tables (CoCoPUTs): Facilitating Genetic Variation
755         Analyses and Recombinant Gene Design. J Mol Biol. 2019 Jun 14;431(13):2434–41.
756   53.   Athey J, Alexaki A, Osipova E, Rostovtsev A, Santana-Quintero LV, Katneni U, et al. A
757         new and updated resource for codon usage tables. BMC Bioinformatics. 2017 Sep
758         2;18(1):391.
759   54.   Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The
760         Molecular Signatures Database Hallmark Gene Set Collection. Cell Syst. 2015 Dec
761         23;1(6):417–25.
762   55.   Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
763         features. Bioinforma Oxf Engl. 2010 Mar 15;26(6):841–2.
764

765 **ADDITIONAL FILES**

766 **Additional file 1 - Supplemental Methods**

767 Format: PDF (.pdf)

768 List of reagents and detailed description the computational software used.

769 **Additional file 2 - Supplemental Figures**

770 Format: PDF (.pdf)

771 **Figure S1.** Medians of Relative Translation Efficiencies weights (RTEw) across all TCGA

772 tissues, related to Figure 3. **Figure S2.** Differential expression of tRNAs between healthy

773 and tumor samples across 22 cancer types, related to Figure 4. **Figure S3.** Differential

774 methylation and copy number between healthy and tumor samples of tRNA genes, related to

775 Figure 5.

776 **Additional file 3 - Table S1. Correlation Matrix**

777 Format: Comma Separated Values (.csv)

778 Correlations of tRNA expression of HEK293 from four small RNA-seq datasets against three

779 conventional tRNA-seq quantifications.

780 **Additional file 4 - Table S2. Sequencing Quality**

781 Format: XLSX (.xlsx)

782 Comparison of tRNA reads coverage between small RNA-seq and hydro-tRNAseq datasets,

783 as well as the coefficient of variance and the standard deviation between replicates.

784 **Additional file 5 - Table S3. TCGA samples**

785 Format: Tab Separated Values (.tsv)

786 Number and abbreviations of TCGA samples covering 23 cancer types.

787    **Additional file 6 - Table S4. GSEA RAA**

788    Format: XLSX (.xlsx)

789    Component features of the Principal Component Analysis of the Relative Anticodon

790    Abundances (RAA, see Figure 2B). GSEA of the correlations of the first PCA component

791    against all genes.

792    **Additional file 7 - Table S5. GSEA RTE**

793    Format: XLSX (.xlsx)

794    Component features of the Principal Component Analysis of the Relative Tranlational

795    Efficiency weights (RTEw, see Figure 3B). GSEA of the correlations of the first two PCA

796    components against all genes.

797    **Additional file 8 - Table S6. GSEA deltaRTE**

798    Format: XLSX (.xlsx)

799    GSEA of the differential RTE between extreme tissues ($\Delta RTE = RTE_{Colorectal} - RTE_{Brain}$).

800    **Additional file 9 - Table S7. RTE survival analysis**

801    Format: Comma Separated Values (.csv)

802    Association of RTEw with cancer prognosis across 22 cancer types. The survival analysis

803    was performed for all codons whose translational efficiency was significantly different in

804    more than 5 cancer types in the one direction with respect to the other [Abs(UP-DOWN)>5],

805    and correspondingly corrected for multiple comparisons using FDR.

806    **Additional file 10 - Table S8. GSEA deltaRTE KIRC**

807    Format: XLSX (.xlsx)

808    GSEA of the differential RTE between cancer and tumor samples from KIRC ($\Delta$RTE =

809    $RTE_{Tumor}$ - $RTE_{Healthy}$).

810    **Additional file 11 - Table S9. GSEA RCU ProCCA**
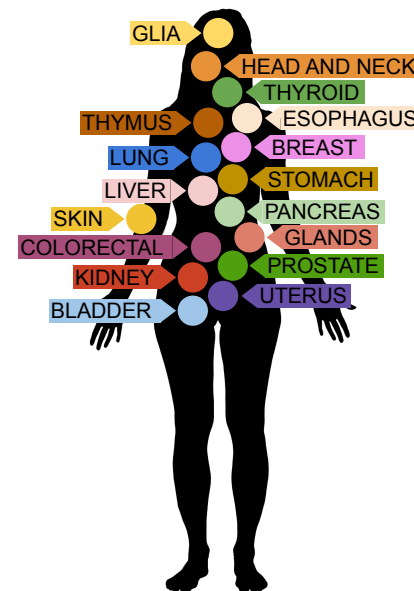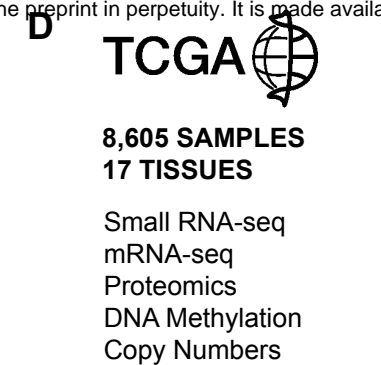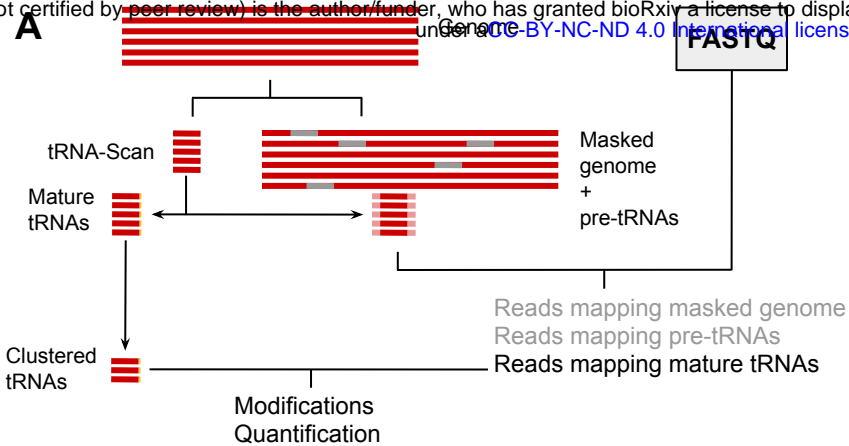
811    Format: XLSX (.xlsx)

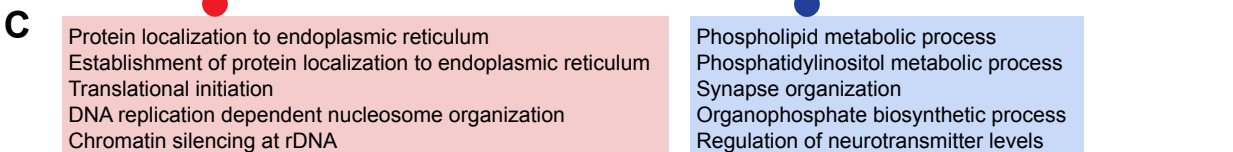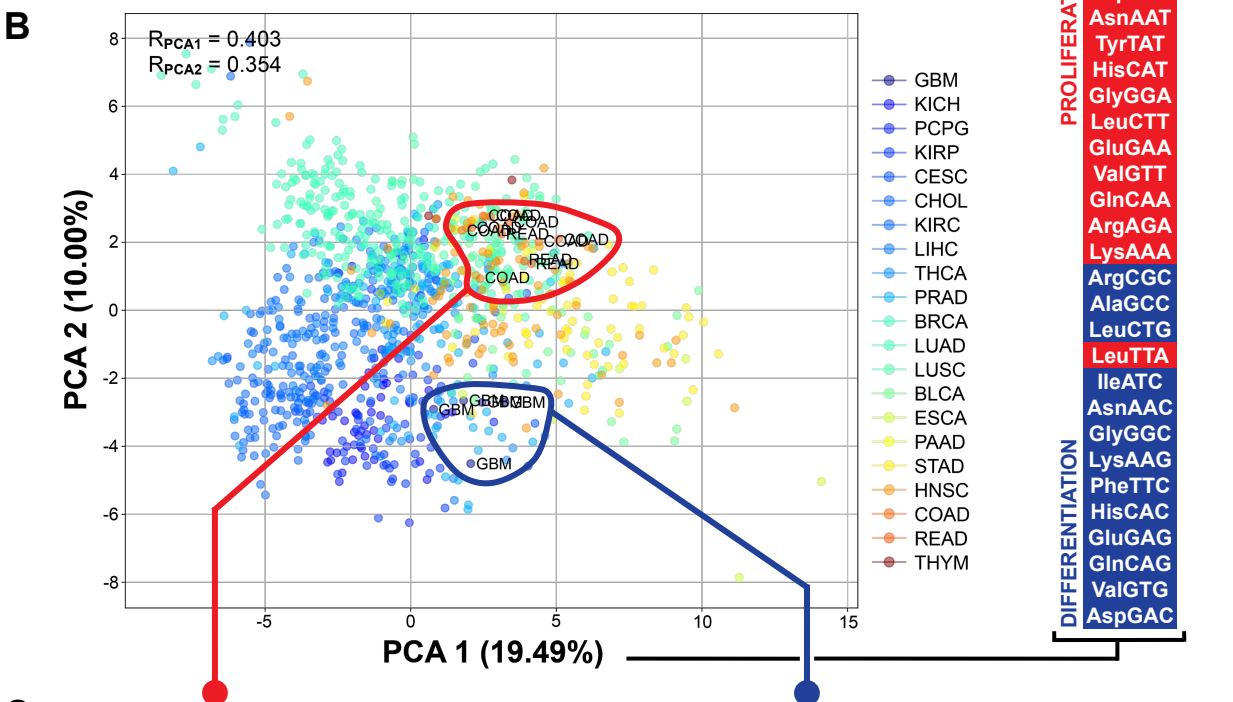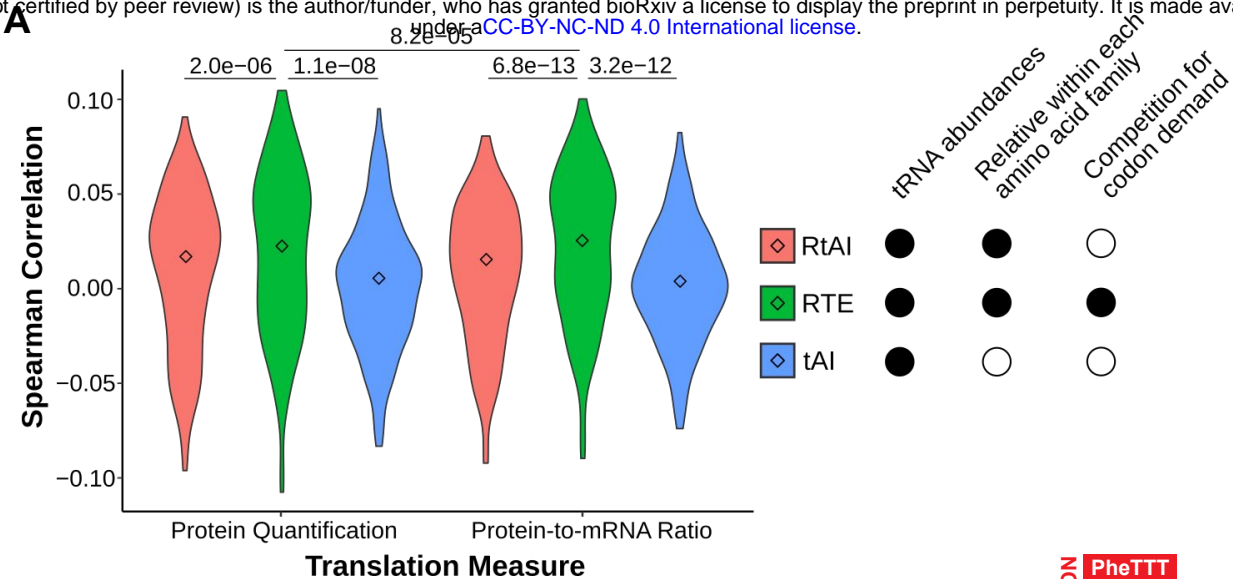812    GSEA of the Relative Codon Usage (RCU) of the codon ProCCA among the whole genome.

813    **Additional file 12 - Table S10. MLR coefficients**

814    Format: XLSX (.xlsx)

815    Multiple Linear Regression (MLR) between the square-root-normalized tRNA expression and
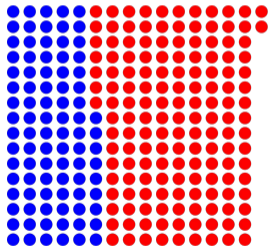
816    the promoter methylation and gene copy number.

817

818

**A**



**B**



**C**

**A**

$$EXPRESSION = \beta_0 + \beta_{Me}Me + \beta_{CNA}CNA$$

$\beta_{Me} < 0$ (p = 4.67e-05)    $\beta_{CNA} > 0$ (p = 1.80e-15)

**B**